

View-based query processing: On the relationship between rewriting, answering and losslessness

Diego Calvanese^{a,*}, Giuseppe De Giacomo^b, Maurizio Lenzerini^b, Moshe Y. Vardi^c

^a Faculty of Computer Science, Free University of Bozen-Bolzano, Piazza Domenicani 3, 39100 Bolzano, Italy

^b Dipartimento di Informatica e Sistemistica, Università di Roma “La Sapienza”, Via Salaria 113, 00198 Roma, Italy

^c Department of Computer Science, Rice University, P.O. Box 1892, Houston, TX 77251-1892, USA

Abstract

As a result of the extensive research in view-based query processing, three notions have been identified as fundamental, namely rewriting, answering, and losslessness. *Answering* amounts to computing the tuples satisfying the query in all databases consistent with the views. *Rewriting* consists in first reformulating the query in terms of the views and then evaluating the rewriting over the view extensions. *Losslessness* holds if we can answer the query by solely relying on the content of the views. While the mutual relationship between these three notions is easy to identify in the case of conjunctive queries, the terrain of notions gets considerably more complicated going beyond such a query class. In this paper, we revisit the notions of answering, rewriting, and losslessness and clarify their relationship in the setting of semistructured databases, and in particular for the basic query class in this setting, i.e., two-way regular path queries. Our first result is a clean explanation of the relationship between answering and rewriting, in which we characterize rewriting as a “linear approximation” of query answering. We show that applying this linear approximation to the constraint-satisfaction framework yields an elegant automata-theoretic approach to query rewriting. As for losslessness, we show that there are indeed two distinct interpretations for this notion, namely with respect to answering, and with respect to rewriting. We also show that the constraint-theoretic approach and the automata-theoretic approach can be combined to give algorithmic characterization of the various facets of losslessness. Finally, we deal with the problem of coping with loss, by considering mechanisms aimed at explaining lossiness to the user.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Query containment; Query rewriting; Query answering; Losslessness; Conjunctive queries; Regular path queries; Semistructured data

1. Introduction

View-based query processing is the problem of computing the answer to a query based on a set of views [1–3]. This problem has recently received much attention in several application areas, such as mobile computing, query

* Corresponding author. Tel.: +39 0471 016160; fax: +39 0471 016009.

E-mail addresses: calvanese@inf.unibz.it (D. Calvanese), degiamaco@dis.uniroma1.it (G. De Giacomo), lenzerini@dis.uniroma1.it (M. Lenzerini), vardi@cs.rice.edu (M.Y. Vardi).

URLs: <http://www.inf.unibz.it/~calvanese/> (D. Calvanese), <http://www.dis.uniroma1.it/~degiamaco/> (G. De Giacomo), <http://www.dis.uniroma1.it/~lenzerini/> (M. Lenzerini), <http://www.cs.rice.edu/~vardi/> (M.Y. Vardi).

optimization, data warehousing, and data integration. A large number of results have been reported in the last years, and several methods have been proposed (see [4] for a recent survey).

As a result of the extensive research in this area, there is proliferation of notions whose relationship to each other is not clear. Fundamentally, there seem to be two basic approaches to view-based query processing. The first approach, originating with [1], is the *query-rewriting* approach, which is based on the idea of first reformulating the query in terms of the views and then evaluating the rewriting over the view extensions. The other approach, originating with [5], is the *query-answering* approach, which takes a more direct route, trying to compute the so-called *certain tuples*, i.e., the tuples satisfying the query in all databases consistent with the views, on the basis of the view definitions and the view extensions. The relationship between the two approaches has been discussed (e.g., [6,7]), but not completely clarified, and is often ignored, see for example [1,8,9].

A related issue that has been studied in several papers is whether the information content of the views is sufficient to answer a given query *completely*, i.e., as if one could access the underlying database. We say that a set of views is *lossless* with respect to a query, if, no matter what the database is, we can answer completely the query by solely relying on the content of the views. This concept has several applications, for example, in view selection [10], where we have to measure the quality of the choice of the views to materialize in the data warehouse, or in data integration, where we may be interested in checking whether the relevant queries can be answered by accessing only a given set of sources [11]. Several papers have addressed the issue of losslessness implicitly [1,12,11] or explicitly [13]. It should be noted, however, that losslessness is relative to the manner in which view-based query processing is performed, since the goal is lossless query processing. Thus, there ought to be two distinct notions of losslessness, with respect to query rewriting or with respect to query answering. Recent discussions of losslessness, such as [11,13], do not reflect this distinction.

One of the reasons the distinction between query answering and query rewriting has been blurred is that much of the work in this area has focused on using conjunctive queries for both target queries and view definitions, cf. [4]. This setting turns out to be extremely well behaved. In particular, query rewriting and query answering coincide, if we allow the target query to be written as a union of conjunctive queries. Furthermore, losslessness with respect to query rewriting and with respect to query answering also coincide, even if we require rewriting by conjunctive queries (disallowing unions). These results, implicit or explicit in [1], give the impression of a simple “terrain” of notions. Once, however, one goes even slightly beyond conjunctive queries or slightly modifies the view model, the terrain of notions gets considerably more complicated, as has already been observed in [3].

In this paper, we revisit the notions of query answering, query rewriting and losslessness, and clarify their relationship in the setting of semistructured databases, which capture data that do not fit into rigid, predefined schemas, and are best described by graph-based data models [14–17]. The prevalent model for semistructured data is that of edge-labeled graphs, in which nodes describe data elements and edges describe relationships or values. (Extensions to node-labeled graphs or to node-edge-labeled graphs are straightforward.)

Methods for extracting information from semistructured data necessarily incorporate special querying mechanisms that are not common in traditional database systems. One such basic mechanism is that of *regular-path queries* (RPQs), which retrieves all pairs of nodes in the graph connected by a path conforming to a regular expression [18,19]. We allow in our regular path queries also the inverse operator. The inverse operator is essential for expressing navigations in the database that traverse the edges both backward and forward [20]. We call such queries *two-way regular path queries* (2RPQs). Such path queries are useful in real settings (see for example [14,18,21]), and are part of the core of many query languages for semistructured data [19,22,23]. In our earlier work we studied both query answering and query rewriting for 2RPQs [24]. For an introductory survey on 2RPQs, see [25].

Our first result is a clean explanation of the relationship between query rewriting and query answering. We view query answering as the more robust notion among the two, since its definition is in terms of the information content of the view extensions. The certain tuples are the tuples whose presence in the answer logically follows from the view extension. In contrast, query rewriting is motivated by the pragmatic need to access the view extensions using a query language that is close, if not identical, to the language in which the target query and the views were formulated. For example, [1] considered rewriting of conjunctive queries by means of unions of conjunctive queries, [26] considered rewriting of RPQs by means of RPQs, and [24] considered rewriting of 2RPQs using 2RPQs.

The setup we use in this paper is that of *sound views*, in which view extension need not reflect global data completely. Thus, all we require from a view V_i defined in terms of a query Q_i is that its extension E_i with respect to a global database \mathcal{B} is such that $E_i \subseteq Q_i(\mathcal{B})$. This setting corresponds to the long-standing *open-world approach*

for querying incomplete information [27]. In this setting query answering can be characterized in terms of constraint satisfaction (or, equivalently, the homomorphism problem [28]), with a constraint template derived from the target query and view definition [7].

A second contribution is the introduction of the notion of “linear approximation” of query answering, and a characterization of the rewriting of a 2RPQ by means of 2RPQs in terms of such a linear approximation. The linear approximation of query answering consists in retrieving a pair (c, d) from the view extension only if its inclusion in the answer is logically implied by a single path between c and d in the view extension. (For 2RPQs two-way paths are considered, while for RPQs one-way paths are considered.) We show that, by applying this linear approximation, we can exploit the constraint-satisfaction framework also for query rewriting. In particular, this yields a natural extension to 2RPQs of the elegant automata-theoretic approach to query rewriting of [26].

As a third contribution, we show that there are indeed two distinct notions of losslessness. Losslessness with respect to query rewriting is what has been called *exactness* in [26], while losslessness with respect to query answering is what has been called simply *losslessness* in [13], and which we view as the more fundamental notion. Since query rewriting is an approximation of query answering, exactness is a stronger notion than losslessness; exactness implies losslessness, but not vice versa. Exactness was taken in [26] to be a measure of quality of query rewriting, but we now see that it combines query rewriting with losslessness. A better way to measure the quality of query rewriting is to measure its quality as an approximation. We say that query rewriting is *perfect* if it is equivalent to query answering. Thus, exactness is the conjunction of perfectness and losslessness (with respect to query answering). We show that the constraint-theoretic approach and the automata-theoretic approach can be combined to give algorithmic characterization of the three notions of perfectness, losslessness, and exactness.

We also consider lossiness, which we view as the central challenge of view-based query processing, as lossiness is more likely to be the norm rather than the exception. Once a schema designer has learned that a view decomposition is lossy with respect to a certain query, how should this “loss” be dealt with? We believe that database design tools should help users to “cope with loss”. In particular, we believe that it would be useful to the user to understand what information is lost by view-based query answering. We discuss a variety of mechanisms aimed at explaining such lossiness to the user.

Finally, we discuss how exactness, perfectness, losslessness, and lossiness relate to each other. In this way, we get a complete picture of the relationships among maximal rewriting, linear approximation, certain answers, and the query itself.

The paper is organized as follows. In Section 2 we recall the basic notions related to view-based query processing, and in Section 3 we recall the relationship between query answering and constraint satisfaction. In Section 4 we discuss the relationship between answering and rewriting. In Section 5 we study losslessness with respect to rewriting for 2RPQs and in Section 6 losslessness with respect to answering. For the latter we introduce the notion of linear fragments of certain answers. In Section 7 we conclude the paper with a final discussion relating the various notions to each other.

2. Preliminaries

Following the usual approach in semistructured data [17], we define a *semistructured database* as a finite directed graph whose edges are labeled by elements from a given finite alphabet Σ . Each node represents an object and an edge from object x to object y labeled by r , denoted $r(x, y)$, represents the fact that relation r holds between x and y . Observe that a semistructured database can be seen as a (finite) relational structure over the set Σ of binary relational symbols. A *relational structure* (or simply *structure*) \mathcal{B} over Σ is a pair $(\Delta^{\mathcal{B}}, \cdot^{\mathcal{B}})$, where $\Delta^{\mathcal{B}}$ is a finite domain and $\cdot^{\mathcal{B}}$ is a function that assigns to each relation symbol in $r \in \Sigma$ a binary relation $r^{\mathcal{B}}$ over $\Delta^{\mathcal{B}}$, also denoted by $r(\mathcal{B})$.

A query is a function from relational structures to relations, assigning to each relational structure over a given alphabet a relation of a certain arity. In this paper we deal mainly with binary queries. A *regular-path query* (RPQ) over Σ is defined in terms of a regular language over Σ . The *answer* $Q(\mathcal{B})$ to an RPQ Q over a database \mathcal{B} is the set of pairs of objects connected in \mathcal{B} by a directed path traversing a sequence of edges forming a word in the regular language $L(Q)$ defined by Q .

RPQs allow for navigating the edges of a semistructured databases only in the forward direction. RPQs extended with the ability of navigating database edges backward are called *two-way regular-path queries* (2RPQs) [24].

Formally, we consider an alphabet $\Sigma^\pm = \Sigma \cup \{r^- \mid r \in \Sigma\}$ which includes a new symbol r^- for each relation symbol r in Σ . The symbol r^- denotes the *inverse* of the binary relation r . If $p \in \Sigma^\pm$, then we use p^- to mean the inverse of p , i.e., if p is r , then p^- is r^- , and if p is r^- , then p^- is r . A 2RPQ over Σ is defined in terms of a regular language over Σ^\pm . The *answer* $Q(\mathcal{B})$ to a 2RPQ Q over a database \mathcal{B} is the set of pairs of objects connected in \mathcal{B} by a semipath that conforms to the regular language $L(Q)$. A *semipath* in \mathcal{B} from x to y (labeled with $p_1 \cdots p_n$) is a sequence of the form $(x_0, p_1, x_1, \dots, x_{n-1}, p_n, x_n)$, where $n \geq 0$, $x_0 = x$, $x_n = y$, and for each x_{i-1}, p_i, x_i , we have that $p_i \in \Sigma^\pm$, and, if $p_i = r$ then $(x_{i-1}, x_i) \in r(\mathcal{B})$, and if $p_i = r^-$ then $(x_i, x_{i-1}) \in r(\mathcal{B})$. Intuitively, a semipath $(x_0, p_1, x_1, \dots, x_{n-1}, p_n, x_n)$ corresponds to a navigation of the database from x_0 to x_n , following edges forward or backward, according to the sequence of edge labels $p_1 \cdots p_n$. Note that the objects in a semipath are not necessarily distinct. A semipath is said to be *simple* if no object in it appears more than once. A *linear database* with endpoints x and y is a database constituted by a single simple semipath from x to y . We say that a semipath $(x_0, p_1, \dots, p_n, x_n)$ *conforms to* a 2RPQ Q if $p_1 \cdots p_n \in L(Q)$. Summing up, a pair (x, y) of objects is in the answer $Q(\mathcal{B})$ if and only if, by starting from x , it is possible to reach y by navigating on \mathcal{B} according to one of the words in $L(Q)$. The notions above can be extended to *two-way path queries*, which are defined similarly to 2RPQs, but without requiring the language to be regular.

Consider now a semistructured database that is accessible only through a collection of views expressed as 2RPQs, and suppose we need to answer a 2RPQ over the database only on the basis of our knowledge on the views. Specifically, the collection of views is represented by a finite set \mathcal{V} of *view symbols*, each denoting a binary relation. Each view symbol $V \in \mathcal{V}$ has an associated *view definition* V^Σ , which is a 2RPQ over Σ . A \mathcal{V} -*extension* \mathcal{E} is a relational structure over \mathcal{V} . We consider views to be *sound* [3,29], i.e., we model a situation where the extension of the views provides a subset of the results of applying the view definitions to the database. Formally, given a set \mathcal{V} of views and a database \mathcal{B} , we use $\mathcal{V}^\Sigma(\mathcal{B})$ to denote the \mathcal{V} -extension \mathcal{E} such that $V(\mathcal{E}) = V^\Sigma(\mathcal{B})$, for each $V \in \mathcal{V}$. We say that a \mathcal{V} -extension \mathcal{E} is *sound wrt a database* \mathcal{B} if $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$. In other words, for a view $V \in \mathcal{V}$, all the tuples in $V(\mathcal{E})$ must appear in $V^\Sigma(\mathcal{B})$, but $V^\Sigma(\mathcal{B})$ may contain tuples not in $V(\mathcal{E})$.

Given a set \mathcal{V} of views, a \mathcal{V} -extension \mathcal{E} , and a query Q over Σ , the set of *certain answers* (under sound views) to Q with respect to \mathcal{V} and \mathcal{E} is the set of pairs (x, y) of objects such that $(x, y) \in Q(\mathcal{B})$ for every database \mathcal{B} wrt which \mathcal{E} is sound, i.e., $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$. *View-based query answering* consists in deciding whether a given pair of objects is a certain answer to Q with respect to \mathcal{V} and \mathcal{E} . Given a set \mathcal{V} of views and a query Q , we denote by $\text{cert}_{Q,\mathcal{V}}$ the query that, for every \mathcal{V} -extension \mathcal{E} , returns the set of certain answers to Q with respect to \mathcal{V} and \mathcal{E} .

View-based query answering has also been tackled using an indirect approach, called *view-based query rewriting*. According to such an approach, a query Q over the database alphabet is processed by first reformulating Q into an expression of a fixed query language over the view alphabet \mathcal{V} (called *rewriting*), and then evaluating the rewriting over the view extensions. Formally, let Q be a query over the database alphabet, and let Q_r be a query over the view alphabet \mathcal{V} . We say that Q_r is a *rewriting of Q under sound views* \mathcal{V} (or simply, with respect to views \mathcal{V}), if for every database \mathcal{B} and for every \mathcal{V} -extension \mathcal{E} with $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$, we have that $Q_r(\mathcal{E}) \subseteq Q(\mathcal{B})$. Since 2RPQs are monotone, by results in [7] (Proposition 13 and 24), rewritings admit the following simpler characterization. A 2RPQ Q_r is a rewriting of a 2RPQ Q if, for every database \mathcal{B} , we have that $Q_r(\mathcal{V}^\Sigma(\mathcal{B})) \subseteq Q(\mathcal{B})$. We make use of this characterization in the following.

Obviously, in view-based query rewriting, we are not interested in arbitrary rewritings, but we aim at computing rewritings that capture the original query at best. Let \mathcal{C} be a query class in which rewritings are expressed. A query Q_r in \mathcal{C} is a \mathcal{C} -*maximal rewriting* of Q under \mathcal{V} if (i) it is a rewriting of Q under \mathcal{V} , and (ii) for each query Q'_r in \mathcal{C} that is a rewriting of Q under \mathcal{V} and for each database \mathcal{B} and each \mathcal{V} -extension \mathcal{E} with $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$, we have that $Q'_r(\mathcal{E}) \subseteq Q_r(\mathcal{E})$.¹ Since in this paper we are focusing on 2RPQs, we are interested in the case where also rewritings are 2RPQs over the view alphabet \mathcal{V} , i.e., rewritings are expressed in the same language as queries over the database.

Throughout the paper, we will assume that RPQs are expressed as finite state automata over an appropriate alphabet. Besides standard (one-way) deterministic and non-deterministic finite state automata over words (1DFAs and 1NFAs, respectively), we assume familiarity with two-way automata (2NFAs) [30].

¹ Observe that, by definition, all maximal rewritings are semantically equivalent, though they may be syntactically different. Hence, with some abuse of terminology, we will talk about “the” maximal rewriting.

3. Answering and constraint satisfaction

In this work we make use of the tight relationship between view-based query answering for RPQs and 2RPQs and constraint satisfaction, which we recall here.

A *constraint-satisfaction problem (CSP)* is traditionally defined in terms of a set of variables, a set of values, and a set of constraints, and asks whether there is an assignment of the variables with the values that satisfies the constraints. A characterization of CSP can be given in terms of homomorphisms between relational structures [28]. Here we consider relational structures whose relations are of arbitrary arity.

A *homomorphism* $h : A \rightarrow B$ between two relational structures A and B over the same alphabet is a mapping $h : \Delta^A \rightarrow \Delta^B$ such that, if $(c_1, \dots, c_n) \in r(A)$, then $(h(c_1), \dots, h(c_n)) \in r(B)$, for every relation symbol r in the alphabet. Let \mathcal{A} and \mathcal{B} be two classes of structures. The *(uniform) constraint-satisfaction problem* $CSP(\mathcal{A}, \mathcal{B})$ is the following decision problem: given a structure $A \in \mathcal{A}$ and a structure $B \in \mathcal{B}$ over the same alphabet, is there a homomorphism $h : A \rightarrow B$? When \mathcal{B} consists of a single structure B and \mathcal{A} is the set of all structures over the alphabet of B , we get the so-called *non-uniform* constraint-satisfaction problem, denoted by $CSP(B)$, where B is fixed and the input is just a structure $A \in \mathcal{A}$. As usual, we use $CSP(B)$ also to denote the set of structures A such that there is a homomorphism from A to B . From the very definition of CSP it follows directly that every $CSP(\mathcal{A}, \mathcal{B})$ problem is in NP.

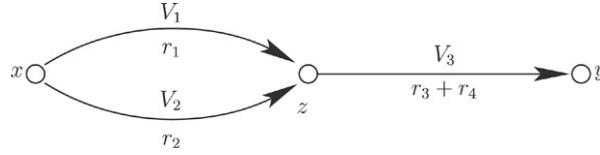
A tight relationship between non-uniform CSP and view-based query answering for RPQs and 2RPQs has been developed in [31,7]. Such a relationship is based on the notions of constraint templates, associated to the query and view definitions, and constraint instance, associated to the view extension. Formally, given a 2RPQ Q and a set \mathcal{V} of 2RPQ views, the *constraint template* $CT_{Q,\mathcal{V}}$ of Q with respect to \mathcal{V} is the relational structure C defined as follows.

- The alphabet of C is $\mathcal{V} \cup \{U_i, U_f\}$, where each view denotes a binary relation symbol, and U_i and U_f are unary relation symbols.
- Let $A^Q = (\Sigma^\pm, S^Q, S_0^Q, \varrho^Q, F^Q)$ be a 1NFA for Q , where Σ^\pm is the alphabet, S^Q is the set of states, S_0^Q is the set of initial states, ϱ^Q is the transition relation, and F^Q is the set of final states. The structure $C = (\Delta^C, \cdot^C)$ is given by:
 - . $\Delta^C = 2^{S^Q}$;
 - . $\sigma \in U_i(C)$ iff $S_0^Q \subseteq \sigma$;
 - . $\sigma \in U_f(C)$ iff $\sigma \cap F^Q = \emptyset$;
 - . for a view $V \in \mathcal{V}$, we have that $(\sigma_1, \sigma_2) \in V^C$ iff there exists a word $q_1 \cdots q_k \in L(V^\Sigma)$ and a sequence T_0, \dots, T_k of subsets of S^Q such that the following hold:
 - (1) $T_0 = \sigma_1$ and $T_k = \sigma_2$,
 - (2) if $s \in T_i$ and $(s, q_{i+1}, t) \in \varrho^Q$ then $t \in T_{i+1}$, for $0 \leq i < k$, and
 - (3) if $s \in T_i$ and $(s, q_i^-, t) \in \varrho^Q$ then $t \in T_{i-1}$, for $0 < i \leq k$.

Intuitively, the constraint template represents for each view V how the states of A^Q (i.e., of the 1NFA for Q) change when we follow database edges according to (paths specified by) words in $L(V^\Sigma)$. Specifically, the last condition above corresponds to saying that a pair of sets of states (σ_1, σ_2) is in $V(C)$ if and only if there is some word w in $L(V^\Sigma)$ such that the following holds: if we start from a state in σ_1 on the left edge of w and move back and forth on w according to the transitions in A^Q , then, if we end up at the left edge of w we can be only in states in σ_1 , and if we end up at the right edge of w we can be only in states in σ_2 ; similarly, if we start from a state in σ_2 on the right edge of w . Moreover, the sets of states in $U_i(C)$ contain all initial states of A^Q , while the sets of states in $U_f(C)$ do not contain any final state of A^Q . This takes into account that we aim at characterizing counterexamples to view-based query answering, and hence we are interested in not getting to a final state of A^Q , regardless of the initial state from which we start and how we follow transitions.

Observe that, to check the existence of a word $q_1 \cdots q_k \in L(V^\Sigma)$ and of a sequence T_0, \dots, T_k of subsets of S such that conditions (1)–(3) above are satisfied, we can resort to a construction analogous to the one in [32]. Hence, such a check can be done in polynomial space in the size of Q , and in fact in nondeterministic logarithmic space in the size of V^Σ .

Given a \mathcal{V} -extension \mathcal{E} and a pair of objects c, d , the *constraint instance* $\mathcal{E}^{c,d}$ is the structure $I = (\Delta^I, \cdot^I)$ over the alphabet $\mathcal{V} \cup \{U_i, U_f\}$ defined as follows:

Fig. 1. The structure of the $\text{cert}_{Q,V}$ in Example 4.1.

- $\Delta^I = \Delta^E \cup \{c, d\}$;
- $V(I) = V(E)$, for each $V \in \mathcal{V}$;
- $U_i(I) = \{c\}$, and $U_f(I) = \{d\}$.

The following theorem provides the characterization of view-based query answering in terms of CSP.

Theorem 3.1 ([7]). *Let Q be a 2RPQ, \mathcal{V} a set of 2RPQ views, \mathcal{E} a \mathcal{V} -extension, and c, d a pair of objects. Then, $(c, d) \notin \text{cert}_{Q,V}(\mathcal{E})$ if and only if there is a homomorphism from $\mathcal{E}^{c,d}$ to $CT_{Q,V}$.*

4. Relationship between rewriting and answering

The relationship between answering and rewriting in view-based query processing is not always well understood. As we said before, one reason for the confusion is that much of the work in this area has focused on a setting based on conjunctive queries, where answering and rewriting coincide. Indeed, if we allow the target query to be written as a union of conjunctive queries (UCQs), then the UCQ-maximal rewriting of the query computes exactly the certain answers. Things get more complicated with RPQs and 2RPQs [6,31], as shown by the following example.

Example 4.1. Consider the RPQ $Q = r_1 \cdot r_3 + r_2 \cdot r_4$ and the views $\mathcal{V} = \{V_1, V_2, V_3\}$ with definitions

$$V_1^\Sigma = r_1, \quad V_2^\Sigma = r_2, \quad V_3^\Sigma = r_3 + r_4.$$

It can be checked that the RPQ-maximal rewriting of Q under \mathcal{V} is empty. On the other hand, $\text{cert}_{Q,V}$ can be expressed as the following conjunctive RPQ [31]

$$\text{cert}_{Q,V} = \{ (x, y) \mid \exists z. x V_1 z \wedge x V_2 z \wedge z V_3 y \}.$$

Note that $\text{cert}_{Q,V}$ matches with non-linear patterns in a database, as depicted in Fig. 1. ■

Interestingly, we show next that we can use the above characterization of view-based query answering in terms of CSP, to characterize also query rewriting, thus providing a clean explanation of the relationship between answering and rewriting.

A preliminary observation is that one can restrict attention to linear databases when looking for counterexamples to rewritings.

Lemma 4.2 ([24]). *Let Q be a 2RPQ, \mathcal{V} a set of 2RPQ views, and w a word over \mathcal{V}^\pm . Then w is not a rewriting (note that w can be viewed as a 2RPQ) of Q with respect to \mathcal{V} if and only if there exists a linear database \mathcal{B} with endpoints c and d , and a view extension \mathcal{E} with $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$, such that $(c, d) \in w(\mathcal{E})$ but $(c, d) \notin Q(\mathcal{B})$.*

Making use of this result, we are able to exploit the constraint template itself as a 1NFA that recognizes the words that do not belong to a rewriting. However, we have first to take care of the fact that only direct view symbols appear in the constraint template, while a rewriting is a 1NFA over direct and inverse view symbols. To do so, we extend the constraint template by adding to the alphabet, for each symbol $V \in \mathcal{V}$, also the inverse symbol V^- . Then we define $(\sigma_1, \sigma_2) \in V^{-C}$ if and only if $(\sigma_2, \sigma_1) \in V^C$. We denote the resulting constraint template with $CT_{Q,V}^\pm$. Observe that the construction of $CT_{Q,V}^\pm$ from $CT_{Q,V}$ takes into account the perfect symmetry that we have when moving along direct and inverse database and view symbols.

Now, $C = CT_{Q,V}^\pm$ can be viewed directly as a 1NFA A^nr over \mathcal{V}^\pm , by taking the domain of C as the set of states of A^nr , the extension of U_i and U_f in C respectively as the set of initial and final states, and by deriving the transition relation of A^nr from the extension of the various $v \in \mathcal{V}^\pm$ as follows: A^nr has a transition (σ_1, v, σ_2) if and only if $(\sigma_1, \sigma_2) \in v^C$.

Let A^{rew} be a 1NFA accepting the complement of A^{nr} . Then the following characterization of the 2RPQ-maximal rewriting holds.

Theorem 4.3. *Let Q be a 2RPQ and \mathcal{V} a set of 2RPQ views. Then A^{rew} is the 2RPQ-maximal rewriting of Q with respect to \mathcal{V} .*

Proof. We show that a word w over \mathcal{V}^\pm is in the 2RPQ-maximal rewriting of Q with respect to \mathcal{V} if and only if $w \notin L(A^{\text{nr}})$.

“ \Rightarrow ” Let $w = v_1 \cdots v_k$ be a word over \mathcal{V}^\pm such that $w \in L(A^{\text{nr}})$. We construct a (linear) view extension \mathcal{E}_w as follows: we introduce pairwise disjoint objects $c = a_0, a_1, \dots, a_k = d$, and for each $i \in \{1, \dots, k\}$,

- $(a_{i-1}, a_i) \in V(\mathcal{E}_w)$, if $v_i = V$, and
- $(a_i, a_{i-1}) \in V(\mathcal{E}_w)$, if $v_i = V^-$.

Consider now the constraint instance $I = \mathcal{E}_w^{c,d}$ obtained from \mathcal{E}_w by defining $U_i(I) = \{c\}$ and $U_f(I) = \{d\}$. It is immediate to verify that, by construction, there is a homomorphism from I to the constraint template $CT_{Q,\mathcal{V}}$. Hence, by Theorem 3.1, we have that $(c, d) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{E}_w)$, while obviously $(c, d) \in w(\mathcal{E}_w)$. Thus w is not part of a rewriting.

“ \Leftarrow ” Let $w = v_1 \cdots v_k$ be a word over \mathcal{V}^\pm that is not part of any rewriting. By Lemma 4.2, there exists a linear database \mathcal{B} containing objects c and d , and a view extension \mathcal{E} with $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$, such that $(c, d) \in w(\mathcal{E})$ but $(c, d) \notin Q(\mathcal{B})$. Since $(c, d) \in w(\mathcal{E})$ and $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$, there are objects $c = a_0, a_1, \dots, a_k = d$ in \mathcal{B} and a word $p = q_1 \cdots q_k$, with each q_i a word over Σ^\pm , such that for each $i \in \{1, \dots, k\}$, we have that $(a_{i-1}, a_i) \in q_i(\mathcal{B})$, and

- $q_i \in L(V^\Sigma)$, if $v_i = V$, and
- $q_i \in L((V^\Sigma)^-)$, if $v_i = V^-$.

We construct from \mathcal{B} a constraint instance $I = \mathcal{E}^{c,d}$ as follows: for each $i \in \{1, \dots, k\}$,

- $(a_{i-1}, a_i) \in V(I)$, if $v_i = V$, and
- $(a_i, a_{i-1}) \in V(I)$, if $v_i = V^-$,

and $U_i(I) = \{c\}$ and $U_f(I) = \{d\}$. Since $(c, d) \notin Q(\mathcal{B})$, and $\mathcal{E} \subseteq \mathcal{V}(\mathcal{B})$, we also have that $(c, d) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{E})$. Hence there is a homomorphism h from I to $CT_{Q,\mathcal{V}}$, and therefore, for each $i \in \{1, \dots, k\}$,

- $(h(a_{i-1}), h(a_i)) \in V(CT_{Q,\mathcal{V}})$, if $v_i = V$, and
- $(h(a_i), h(a_{i-1})) \in V(CT_{Q,\mathcal{V}})$, if $v_i = V^-$.

In any case, this means that there are transitions $(h(a_{i-1}), v_i, h(a_i)) \in Q^{\text{nr}}$ that lead from an initial to a final state of A^{nr} . Hence, $w = v_1 \cdots v_k$ is accepted by A^{nr} . \square

The above characterization provides a nice combination of the constraint based [31] and automata theoretic [24] approaches to view-based query processing for 2RPQs, and goes into the heart of view-based rewriting. A (language) rewriting accepts a pair (c, d) if there is a path between c and d such that, if we view this path as a linear view extension, then (c, d) is in the certain answer with respect to this view extension. That means that there is no homomorphism from this path into the constraint template. Indeed, for a path, the existence of a homomorphism into the constraint template means that the path is accepted by the template, viewed as an automaton. Naturally, the difference with view-based query answering, is that we are not limited to linear view extensions only. Suppose that V_i and V_j connect the same pair of objects in a view extension. In rewriting we have to ignore this and allow the choice of distinct pairs of objects for the two views in a counterexample database. Query answering instead takes into account that the two pairs of objects are the same. Thus, query answering is more precise than query rewriting. On the other hand, the simplification introduced by query rewriting allows one to have polynomial time evaluation in the size of the data, while query answering is coNP-complete [6].

Finally, observe that the above construction provides also optimal upper bounds for the problems of computing the 2RPQ-maximal rewriting and of determining whether such a rewriting is nonempty [24]. Indeed, the constraint template, and hence the 1NFA A^{nr} can be constructed in EXPTIME and is of exponential size [7]. Hence, its complement A^{rew} , which provides the 2RPQ-maximal rewriting, is of double exponential size and can be constructed in 2EXPTIME. On the other hand, if we only want to check its emptiness, we can complement A^{nr} on the fly, getting an EXPSPACE upper bound. All these bounds are tight [26].

5. Losslessness with respect to rewriting

We deal now with the problem of analyzing the loss of information in view-based query processing, and of characterizing the quality of certain answers and of rewritings. For this purpose, we make use of the following basic notions.

- To determine whether the information content of a set of views is sufficient to answer completely a given query, we make use of the notion of losslessness [12,13]. In [13], a set of views \mathcal{V} is said to be *lossless* with respect to a query Q , if for every database \mathcal{B} we have that $Q(\mathcal{B}) = \text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B}))$.
- As for rewritings, equivalence of a rewriting to the original query, modulo the view definitions, is called exactness (cf. [1,26]). Formally, a rewriting Q_r in a certain query class \mathcal{C} is an *exact rewriting* of Q with respect to views \mathcal{V} , if for every database \mathcal{B} we have that $Q(\mathcal{B}) = Q_r(\mathcal{V}^\Sigma(\mathcal{B}))$.
- Finally, to determine whether we lose answering power by resorting to rewriting, we can compare rewritings with the certain answers, with the aim of checking whether the two are actually equivalent. A rewriting Q_r in a certain query class \mathcal{C} is a *perfect rewriting* of Q with respect to views \mathcal{V} , if for every database \mathcal{B} and every view extension \mathcal{E} with $\mathcal{E} \subseteq \mathcal{V}^\Sigma(\mathcal{B})$ we have that $\text{cert}_{Q,\mathcal{V}}(\mathcal{E}) = Q_r(\mathcal{E})$.

The first notion aims at determining possible loss with respect to view-based query answering, and will be discussed in the next section. The other two notions deal with the loss of information in the case of rewritings, and are discussed below.

In the case of conjunctive queries, the best rewriting of a conjunctive query Q is a union of conjunctive queries. Therefore, checking exactness amounts to verifying whether Q is contained in the UCQ-maximal rewriting. The latter is a, possibly exponential, union of conjunctive queries, each of linear size. Since a conjunctive query is contained in a union of conjunctive queries only if it is contained in one of its disjuncts, it suffices to check whether there is a single conjunctive query in the rewriting that is equivalent to Q , after substituting the view definitions. This can be done in NP in the size of Q . As for perfectness, we already observed that the maximal rewriting computes exactly the certain answers. Therefore, the maximal rewriting is always perfect.

In the case of 2RPQs, things are more complicated. Exactness is studied in [24], where it is shown that verifying the existence of an exact rewriting is 2EXPSpace-complete. On the other hand, perfectness is a new notion, and we provide here a method for checking perfectness of the 2RPQ-maximal rewriting A^{rew} of a query Q . Exploiting the fact that 2RPQs are monotone, by results in [7], this amounts to check whether for all databases \mathcal{B} we have that $\text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B})) \subseteq A^{\text{rew}}(\mathcal{V}^\Sigma(\mathcal{B}))$. This corresponds to checking whether Q is *view-based contained* in A^{rew} (see [7]). To do this check, we can in principle directly use the technique in [7]. Since the 2RPQ-maximal rewriting A^{rew} is a 1NFA of double exponential size in Q , and checking whether Q is view-based contained in A^{rew} can be done in NEXPTIME in Q and A^{rew} [7], this gives us a N3EXPTIME upper bound. However, we can do better, by making use of the fact that we have obtained the 1NFA A^{rew} for the rewriting by complementation, and thus by application of the subset construction. This allows us to characterize non-membership in the answer set to A^{rew} by homomorphism into a structure $C = (\Delta^C, \cdot^C)$, called the *rewriting constraint template* $\text{CTR}_{A^{\text{rew}},\mathcal{V}}$ of A^{rew} , defined as follows:

- The alphabet of C is $\mathcal{V}^\pm \cup \{U_i, U_f\}$, where U_i and U_f denote unary relation symbols.
- Let $A^{\text{nr}} = (\mathcal{V}^\pm, S, S_0, \varrho, F)$ be a 1NFA for the complement of the rewriting (see Section 4). Then
 - $\Delta^C = 2^S$;
 - $\sigma \in U_i^C$ iff $S_0 \subseteq \sigma$;
 - $\sigma \in U_f^C$ iff $\sigma \subseteq F$;
 - $(\sigma_1, \sigma_2) \in r^C$ iff $\varrho(\sigma_1, r) \subseteq \sigma_2$ and $\varrho(\sigma_2, r^-) \subseteq \sigma_1$.

To characterize perfectness of the rewriting in terms of CSP, we need to introduce proper constraint templates (see also [7]). Given the rewriting constraint template $\text{CT}_{A^{\text{rew}},\mathcal{V}}$, a *proper constraint template* $\text{CT}_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$ is obtained by eliminating from \mathcal{U}_i all but one element α and from \mathcal{U}_f all but one element β .

Lemma 5.1. *Let Q be a 2RPQ and \mathcal{V} be a set of 2RPQ views. Then the 2RPQ-maximal rewriting of Q with respect to \mathcal{V} is perfect if and only if for every proper constraint template $\text{CTR}_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$ of $\text{CTR}_{A^{\text{rew}},\mathcal{V}}$, there exists a homomorphism from $\text{CTR}_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$ to $\text{CT}_{Q,\mathcal{V}}$.*

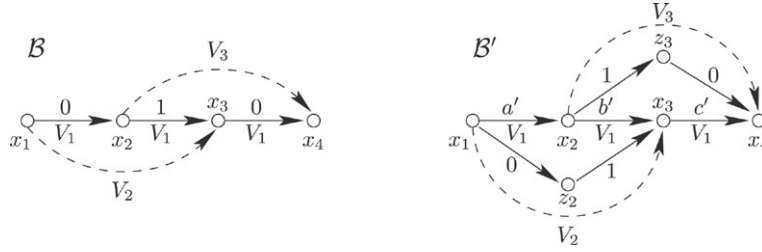


Fig. 2. Databases of Example 6.1.

Proof. Consider now a view extension $\mathcal{E} = \mathcal{V}^\Sigma(\mathcal{B})$ for some database \mathcal{B} , and a pair (c, d) that is not in the answer set $A^{\text{rew}}(\mathcal{E})$. This means that every semipath $(x_0, v_1, x_1, \dots, x_{n-1}, v_n, x_n)$ from $c = x_0$ to $d = x_n$ in \mathcal{E} can be instantiated by words $w_i \in L(v_i)$ such that Q , when evaluated over the linear database $(x_0, w_1, x_1, \dots, x_{n-1}, w_n, x_n)$, does not accept (c, d) . Thus, we can label the objects x_0, \dots, x_n on the path from c to d in the view extension by sets of states of A^{nr} (recall that A^{nr} is the INFA from which we have obtained A^{rew} by complementation) such that the label of c contains the start set and the label of d consists of accepting states of A^{nr} . This holds for all paths. Thus, we should be able to label all the objects in the view extension \mathcal{E} by sets of states of A^{nr} such that the label of c contains the start set, the label of d consists of accepting states of A^{nr} , and the edges respect the transition relation of A . In summary, we have mapped the view extension \mathcal{E} into the rewriting constraint template of A^{rew} . Thus, non-membership in the answer-set $A^{\text{rew}}(\mathcal{E})$ is characterized by homomorphism from $\mathcal{E}^{c,d}$ into the rewriting constraint template of A^{rew} . Now, non-membership in $\text{cert}_{Q,\mathcal{V}}(\mathcal{E})$ is also characterized by homomorphism from $\mathcal{E}^{c,d}$ to $CT_{Q,\mathcal{V}}$. Hence, by reasoning as in the proof of Theorem 18 in [7], we get the claim. \square

The above characterization provides us with a tighter upper bound than the one discussed above.

Theorem 5.2. Let Q be a 2RPQ and \mathcal{V} a set of 2RPQ views. Then checking whether the 2RPQ-maximal rewriting of Q with respect to \mathcal{V} is perfect can be done in N2EXPTIME in the size of Q and in NEXPTIME in the size of \mathcal{V}^Σ .

Proof. We have to check whether, for every proper constraint template $CTR_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$ of $CTR_{A^{\text{rew}},\mathcal{V}}$ there exists a homomorphism to $CT_{Q,\mathcal{V}}$. Constructing the constraint template $CT_{Q,\mathcal{V}}$ is in EXPTIME in the size of Q and polynomial in the size of \mathcal{V}^Σ . Constructing the INFA A^{nr} is in EXPTIME in the size of Q and polynomial in the size of \mathcal{V}^Σ . Hence, constructing the rewriting constraint template $CTR_{A^{\text{rew}},\mathcal{V}}$ of A^{rew} and each proper constraint template $CTR_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$ is in 2EXPTIME in the size of Q and in EXPTIME in the size of \mathcal{V}^Σ . Checking the existence of each homomorphism is NP in the size of the $CTR_{A^{\text{rew}},\mathcal{V}}^{\alpha,\beta}$. Moreover, the number of proper constraint templates is $2 \cdot 2^k$, where k is the number of states of A^{nr} , and hence it is doubly exponential in the size of Q . These bounds give us the claim. \square

We conjecture that such an upper bound is tight.

6. Losslessness with respect to answering

We now turn to verifying losslessness with respect to answering. We want to verify whether a set of views \mathcal{V} is lossless with respect to a query Q , i.e., verifying whether $\text{cert}_{Q,\mathcal{V}}$ is equivalent to Q (cf. [13]).

In the case of conjunctive queries, we already observed that the maximal rewriting computes exactly the certain answers. Therefore, losslessness with respect to answering and losslessness with respect to rewriting coincide. This is not necessarily the case for RPQs and 2RPQs, as shown in the following example.

Example 6.1. Consider the RPQ $Q = 010 + 101 + 000 + 111$ and the views $\mathcal{V} = \{V_1, V_2, V_3, V_4, V_5\}$ with definitions

$$V_1^\Sigma = 0 + 1, \quad V_2^\Sigma = 01, \quad V_3^\Sigma = 10, \quad V_4^\Sigma = 000, \quad V_5^\Sigma = 111.$$

The set \mathcal{V} of views is not lossless with respect to rewriting. Indeed, the maximal RPQ-rewriting of Q under \mathcal{V} is $V_4 + V_5$, and hence is not exact, i.e., equivalent to Q .

However, \mathcal{V} is lossless with respect to answering, i.e., $\text{cert}_{Q,\mathcal{V}}$ is equivalent to Q . To show this, consider a linear database $\mathcal{B} = (x_1, a, x_2, b, x_3, c, x_4)$ such that $Q(\mathcal{B})$ is not empty. We consider the case where $abc = 010$ (see Fig. 2).

The case where $abc = 101$ is symmetric, and the cases where $abc = 000$ and $abc = 111$ are trivial. Then, $V_1(\mathcal{B})$ returns all pairs of nodes, while

$$V_2(\mathcal{B}) = \{(x_1, x_3)\}, \quad V_3(\mathcal{B}) = \{(x_2, x_4)\}, \quad V_4(\mathcal{B}) = \emptyset, \quad V_5(\mathcal{B}) = \emptyset.$$

Consider now a (possibly non-linear) database \mathcal{B}' such that $\mathcal{V}(\mathcal{B}) \subseteq \mathcal{V}(\mathcal{B}')$ but with $(x_1, x_4) \notin Q(\mathcal{B}')$ (see Fig. 2). Then, by $V_1(\mathcal{B}) \subseteq V_1(\mathcal{B}')$, there is a path in \mathcal{B}' of the form $(x_1, a', x_2, b', x_3, c', x_4)$; by $V_2(\mathcal{B}) \subseteq V_2(\mathcal{B}')$, there is a path in \mathcal{B}' of the form $(x_1, 0, z_2, 1, x_3)$; and by $V_3(\mathcal{B}) \subseteq V_3(\mathcal{B}')$, there is a path in \mathcal{B}' of the form $(x_2, 1, z_3, 0, x_4)$. Since $(x_1, x_4) \notin Q(\mathcal{B}')$, the path $(x_1, a', x_2, b', x_3, c', x_4)$ must be such that $a'b'c' \in \{001, 011, 100, 110\}$. But then, if $a' = 0$ we have in \mathcal{B}' the path $(x_1, 0, x_2, 1, z_3, 0, x_4)$, and hence $(x_1, x_4) \in Q(\mathcal{B}')$. Similarly, if $a' = 1$ then $c' = 0$ and hence we have in \mathcal{B}' the path $(x_1, 0, z_2, 1, x_3, 0, x_4)$, and again $(x_1, x_4) \in Q(\mathcal{B}')$. In both cases we get a contradiction. ■

Losslessness with respect to answering for RPQs was studied in [13]. In the rest of this section we study losslessness with respect to answering for 2RPQs.

The main step toward this goal is to characterize the linear fragment of certain answers. Formally, the *linear fragment of certain answers* $\text{clin}_{Q,\mathcal{V}}$ for a 2RPQ Q with respect to a set \mathcal{V} of 2RPQ views is the maximal two-way path query² Q' over Σ such that, for every database \mathcal{B} we have that $Q'(\mathcal{B}) \subseteq \text{cert}_{Q,\mathcal{V}}(\mathcal{V}(\mathcal{B}))$. Q' being *maximal* means that, for all path queries Q'' over Σ satisfying the condition above, we have that $Q''(\mathcal{B}) \subseteq Q'(\mathcal{B})$ for every database \mathcal{B} . The following result shows that, in order to characterize the linear fragment of certain answers it is sufficient to restrict attention to linear databases, i.e., databases constituted by a single semipath.

Lemma 6.2. *Let Q' be a two-way path query. Then, if there is a database \mathcal{B} and a pair of objects (c, d) in \mathcal{B} such that $(c, d) \in Q'(\mathcal{B})$ and $(c, d) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B}))$, then there is a linear database \mathcal{B}_ℓ with endpoints c' and d' such that $(c', d') \in Q'(\mathcal{B}_\ell)$ and $(c', d') \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B}_\ell))$.*

Proof. The proof is analogous to the proof of Theorem 9 in [13]. □

Hence, to construct the linear fragment of certain answers, we characterize the set of linear databases of the form $\mathcal{B} = (x_0, q_1, x_1, q_2, \dots, q_m, x_m)$, for some m , such that $(x_0, x_m) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{V}(\mathcal{B}))$. By Theorem 3.1, this holds if and only if there is a homomorphism from the constraint instance $\mathcal{V}(\mathcal{B})^{x_0, x_m}$ to the constraint template $CT_{Q,\mathcal{V}}$. In other words, $(x_0, x_m) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{V}(\mathcal{B}))$ if and only if there is a function $\ell(\cdot)$ (i.e., the homomorphism) that labels x_0, \dots, x_m with sets of states of the 1NFA $A^Q = (\Sigma^\pm, S^Q, S_0^Q, Q^Q, F^Q)$ for Q such that the following conditions (which we call *CT-conditions*) hold:

- $S_0^Q \subseteq \ell(x_0)$;
- $\ell(x_m) \cap F^Q = \emptyset$;
- for each pair of objects x_j and x_h in \mathcal{B} and each view V in \mathcal{V} , we have that, if $(x_j, x_h) \in V^\Sigma(\mathcal{B})$ then there exists a word $q_1 \dots q_k \in L(V^\Sigma)$ and a sequence T_0, \dots, T_k of subsets of S^Q such that the following hold:
 - (1) $T_0 = \ell(x_j)$ and $T_k = \ell(x_h)$,
 - (2) if $s \in T_i$ and $(s, q_{i+1}, t) \in Q^Q$ then $t \in T_{i+1}$, for $0 \leq i < k$, and
 - (3) if $s \in T_i$ and $(s, q_i^-, t) \in Q^Q$ then $t \in T_{i-1}$, for $0 < i \leq k$.

Thus, we are looking for words of the form $\ell_0, q_1, \dots, q_m, \ell_m$, where each ℓ_i is a set of states of A^Q , representing $\ell(x_i)$, and that satisfies the above conditions. As shown by the following lemma, we can construct a 1NFA that accepts such words, and then project away the ℓ_i transitions.

For a word $w \in \Sigma^{\pm*}$, we denote with $\mathcal{B}_w^{a,b}$ the linear database constituted by a path from a to b spelled by w (with arbitrary intermediate nodes).

Lemma 6.3. *Let Q be a 2RPQ and \mathcal{V} be a set of 2RPQ views. Then we can construct in double exponential time in Q and \mathcal{V}^Σ two INFAs $A^{\text{nl}}_{\text{lin}}$ and A^{lin} such that:*

² Recall from Section 2 that two-way path queries are a generalization of 2RPQs in which the language used to define a query is not required to be regular.

- A^{lin} accepts all words $w \in \Sigma^{\pm*}$ such that $(a, b) \notin \text{cert}_{Q, \mathcal{V}}(\mathcal{V}(\mathcal{B}_w^{a,b}))$.
- A^{lin} accepts all words $w \in \Sigma^{\pm*}$ such that $(a, b) \in \text{cert}_{Q, \mathcal{V}}(\mathcal{V}(\mathcal{B}_w^{a,b}))$.

Proof. We construct a 1DFA $A^{\mathcal{V}}$ over the alphabet $\Sigma^{\pm} \cup \Lambda$, with $\Lambda = 2^{S^Q}$, as follows:

- We construct a 1DFA A^0 accepting words that
 - . start with a symbol $\ell \in \Lambda$ such that $S_0^Q \subseteq \ell$, and
 - . end with a symbol $\ell \in \Lambda$ such that $\ell \cap F^Q = \emptyset$.
 These conditions can be checked by A^0 in a straightforward manner using a constant number of states.
- For each view $V \in \mathcal{V}$, we construct a 1NFA A^V as follows. We first define a binary relation V^r on Λ as follows: $(\ell, \ell') \in V^r$ if there exists a word $q_1 \cdots q_k \in L(V^{\Sigma})$ and a sequence T_0, \dots, T_k of subsets of S^Q such that the following hold:
 - (1) $T_0 = \ell$ and $T_k = \ell'$,
 - (2) if $s \in T_i$ and $(s, q_{i+1}, t) \in \varrho^Q$ then $t \in T_{i+1}$, for $0 \leq i < k$, and
 - (3) if $s \in T_i$ and $(s, q_i^-, t) \in \varrho^Q$ then $t \in T_{i-1}$, for $0 < i \leq k$.

Note that for each ℓ, ℓ' we can decide in PSPACE whether $(\ell, \ell') \in V^r$ holds,³ so all the relations V^r can be constructed in time exponential in the size of A^Q and linear in the size of \mathcal{V}^{Σ} . From the relation V^r we can immediately construct a 1NFA that accepts a word starting with ℓ and ending with ℓ' if and only if $(\ell, \ell') \notin V^r$.

We then construct the 1NFA A^V in such a way that it reads a word w , guesses positions j and h , checks that $(\ell_j, \ell_h) \notin V^r$, and checks that $(\ell_j, \ell_h) \in V^{\Sigma}(\mathcal{B}_w)$, where \mathcal{B}_w is a (linear) database obtained from w by considering each symbol ℓ_i as an object. To do so, we first construct a 2NFA that reads a word w , guesses positions j and h , and checks that $(\ell_j, \ell_h) \in V^{\Sigma}(\mathcal{B}_w)$. Note that the 2NFA does not need to consider the actual symbols ℓ_i appearing in w , but only the symbols q_i , hence it has a number of states that is polynomial in V^{Σ} (and does not depend on Q). Then we transform the 2NFA into a 1NFA and modify it in the following way to obtain A^V : when the 1NFA guesses the position j it also remembers ℓ_j and, when it guesses position h , it checks that $(\ell_j, \ell_h) \notin V^r$. The number of states of the 1NFA derived from the 2NFA is exponential, and the necessity to remember ℓ_j along its computation multiplies such a number by at most another exponential in Q . Hence, the overall number of states of A^V is exponential both in Q and in V .

- Finally, we take the union of the automata A^V , for each $V \in \mathcal{V}$, complement it, and intersect with A^0 , thus getting $A^{\mathcal{V}}$. Note that such an automaton is deterministic, and has a number of states that is doubly exponential in the size of Q and \mathcal{V}^{Σ} .

Since we are interested in words that consist of an alternation of symbols in Λ and symbols in Σ^{\pm} , we also construct a 1DFA A^{alt} accepting such words. The number of states of A^{alt} is constant.

Now, using $A^{\mathcal{V}}$ and A^{alt} we can characterize both the linear fragment of the certain answers and its complement.

- By intersecting $A^{\mathcal{V}}$ with A^{alt} , and projecting away the symbols ℓ_i , we obtain an 1NFA A^{lin} .
- Similarly, by intersecting the complement⁴ of $A^{\mathcal{V}}$ with A^{alt} , and projecting away the symbols ℓ_i , we obtain an 1NFA A^{lin} .

Finally, we prove now that A^{lin} accepts a word $w \in \Sigma^{\pm*}$ if and only if $(a, b) \notin \text{cert}_{Q, \mathcal{V}}(\mathcal{V}(\mathcal{B}_w^{a,b}))$. The proof for A^{lin} is analogous. “ \Leftarrow ” Let $w = q_1 q_2 \cdots q_m$ be a word such that $(x_0, x_m) \notin \text{cert}_{Q, \mathcal{V}}(\mathcal{V}(\mathcal{B}_w^{x_0, x_m}))$, and let $\mathcal{B}_w^{x_0, x_m}$ be $(x_0, q_1, x_1, q_2, \dots, q_m, x_m)$. By Theorem 3.1 there exists a labeling $\ell(\cdot)$ of the objects x_0, \dots, x_m with states of A^Q such that the CT-conditions hold. Now consider the word $\ell(x_0)q_1\ell(x_1)q_2 \cdots q_m\ell(x_m)$. One can check that, by construction, such a word is accepted by $A^{\mathcal{V}}$ and A^{alt} . Hence, w is accepted by A^{lin} .

“ \Rightarrow ” If A^{lin} accepts a word $w = q_1 q_2 \cdots q_m$, then there exists a word $\ell_0 q_1 \ell_1 q_2 \cdots q_m \ell_m$ accepted by $A^{\mathcal{V}}$ and A^{alt} . Now, consider the linear database \mathcal{B} of the form $(x_0, q_1, x_1, q_2, \dots, q_m, x_m)$ and a labeling $\ell(\cdot)$ of the objects x_0, \dots, x_m with states of A^Q defined by $\ell(x_i) = \ell_i$, for $i \in \{0, \dots, m\}$. Since $\ell_0 q_1 \ell_1 q_2 \cdots q_m \ell_m$ is accepted by $A^{\mathcal{V}}$, the CT-conditions hold and hence, by Theorem 3.1 $(x_0, x_m) \notin \text{cert}_{Q, \mathcal{V}}(\mathcal{V}(\mathcal{B}))$. \square

³ The condition can be verified by checking the nonemptiness of a 1NFA built as in [33]. The 1NFA has an exponential number of states, however nonemptiness can be checked on the fly, without actually building the 1NFA.

⁴ Observe that $A^{\mathcal{V}}$ is deterministic, so complementing simply amounts to swapping final and non-final states.

Both 1NFAs $A^{\text{nlín}}$ and A^{lin} have a number of states that is doubly exponential in both Q and \mathcal{V}^Σ . Obviously, the two automata accept complementary languages. However, in the proof of the above lemma we show how to construct A^{lin} directly, instead of complementing $A^{\text{nlín}}$, to avoid an additional exponential blowup.

Theorem 6.4. *Let Q be a 2RPQ and \mathcal{V} be a set of 2RPQ views, and $A^{\text{nlín}}$ and A^{lin} the 1NFAs defined as above. Then A^{lin} is the linear fragment $\text{clin}_{Q,\mathcal{V}}$ of the certain answers of Q with respect to \mathcal{V} .*

Corollary 6.5. *The linear fragment of a 2RPQ with respect to a set of 2RPQ views is a 2RPQ.*

Now we can deal with checking losslessness with respect to answering. To check whether a set \mathcal{V} of 2RPQ views is lossless with respect to a 2RPQ query Q , we have to check whether for all databases \mathcal{B} , we have that $Q(\mathcal{B})$ is contained in the certain answers $\text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B}))$. Since Q is itself a 2RPQ, and hence a path query, it suffices to check whether Q is contained in the linear fragment of the certain answers, i.e., whether for all databases \mathcal{B} we have that $Q(\mathcal{B}) \subseteq \text{clin}_{Q,\mathcal{V}}(\mathcal{B})$. By exploiting the characterization of the linear fragment of the certain answers in terms of 1NFAs provided above, we get the following upper bound, which is tight already for RPQs [13].

Theorem 6.6. *Let Q be a 2RPQ and \mathcal{V} be a set of 2RPQ views. Then checking whether \mathcal{V} is lossless with respect to Q can be done in EXPSpace in the size of Q and \mathcal{V}^Σ .*

Proof. By Theorem 6.4, the linear fragment $\text{clin}_{Q,\mathcal{V}}$ of the certain answers is given by A^{lin} . Its complement $A^{\text{nlín}}$ is an 1NFA with a number of states that is doubly exponential in Q and \mathcal{V}^Σ . Recall that for a word $w \in \Sigma^{\pm*}$, we denote with $\mathcal{B}_w^{a,b}$ the linear database constituted by a path from a to b spelled by w (with arbitrary intermediate nodes), and that $A^{\text{nlín}}$ accepts all words w such that $(a, b) \notin \text{cert}_{Q,\mathcal{V}}(\mathcal{V}(\mathcal{B}_w^{a,b}))$. To check whether Q is not contained in $\text{clin}_{Q,\mathcal{V}}$, it suffices to check whether there is a word $w \in \Sigma^{\pm*}$ such that $(a, b) \in Q(\mathcal{B}_w^{a,b})$ and $w \in L(A^{\text{nlín}})$. By Lemma 2 in [25], it suffices to check whether there is a word $w \in L(Q) \cap L(A^{\text{nlín}})$, i.e., whether the intersection of Q and $A^{\text{nlín}}$ is nonempty. $A^{\text{nlín}}$ has a number of states that is double exponential in Q and \mathcal{V} . Considering that emptiness of 1NFAs is NLOGSPACE in the number of states and that the construction of $A^{\text{nlín}}$ can be done on the fly while checking for emptiness, we get the claim. \square

Observe that when we have that a set of views is lossless with respect to a query, we have also, as a side effect, that the linear fragment of certain answers is equivalent to the certain answers, since both are equivalent to the query. Now it is natural to try to understand when the linear fragment of certain answers is equivalent to the certain answers, independently of losslessness with respect to answering. Indeed, in this case, since the certain answers are actually expressible as a 2RPQ over the database, we directly get a characterization of the certain answers in the same language used for expressing the query and thus in terms that are understandable to the user.

Given a 2RPQ Q and a set of 2RPQ views \mathcal{V} , checking whether the linear fragment of certain answers is equivalent to the certain answers amounts to checking whether for every database \mathcal{B} we have that $\text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B})) \subseteq \text{clin}_{Q,\mathcal{V}}(\mathcal{B})$. Consider the 1NFA A^{lin} , constructed above, recognizing the linear fragment $\text{clin}_{Q,\mathcal{V}}$ of the certain answers of Q . One can verify that the certain answers $\text{cert}_{A^{\text{lin}},\mathcal{V}}$ of A^{lin} with respect to \mathcal{V} are actually equivalent to A^{lin} itself. Hence, the above check amounts to verifying whether for all databases \mathcal{B} , we have that $\text{cert}_{Q,\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B})) \subseteq \text{cert}_{A^{\text{lin}},\mathcal{V}}(\mathcal{V}^\Sigma(\mathcal{B}))$. This is a form of view-based containment, and by [7] it can be done in NEXPTIME in the size of Q and A^{lin} . Considering that A^{lin} has a number of states that is doubly exponential in the size of Q and \mathcal{V}^Σ , we get the following upper bound.

Theorem 6.7. *Let Q be a 2RPQ and \mathcal{V} be a set of 2RPQ views. Then checking whether the certain answers $\text{cert}_{Q,\mathcal{V}}$ of Q with respect to \mathcal{V} is equivalent to its linear fragment can be done in N3EXPTIME in the size of Q and \mathcal{V}^Σ .*

We conjecture that such an upper bound can be improved.

7. Discussion

In this paper, we have revisited the notions of answering, rewriting and losslessness in the context of view-based query processing in semistructured databases. In particular the richness of RPQs and 2RPQs allows us to uncover several subtle distinctions between the notions of rewriting and answering, and losslessness with respect to them. Such distinctions are completely blurred when focusing on conjunctive queries, due to the fact that rewriting and answering collapse.

Let Q be a 2RPQ, \mathcal{V} a set of 2RPQ views, and let $R_{Q,\mathcal{V}}^{\max}$ denote the 2RPQ-maximal rewriting of Q with respect to \mathcal{V} . Then, by definition and by results in [7] exploiting the fact that 2RPQs are monotone, we know that for every database \mathcal{B} , the following holds:

$$R_{Q,\mathcal{V}}^{\max}(\mathcal{V}^{\Sigma}(\mathcal{B})) \subseteq^{(1)} \text{clin}_{Q,\mathcal{V}}(\mathcal{B}) \subseteq^{(2)} \text{cert}_{Q,\mathcal{V}}(\mathcal{V}^{\Sigma}(\mathcal{B})) \subseteq^{(3)} Q(\mathcal{B}).$$

Notice that we start from a database \mathcal{B} and are evaluating $\text{cert}_{Q,\mathcal{V}}$ and $R_{Q,\mathcal{V}}^{\max}$ over a particular view extension, namely $\mathcal{V}^{\Sigma}(\mathcal{B})$, instead of an arbitrary view extension \mathcal{E} that is sound with respect to \mathcal{B} , i.e., such that $\mathcal{E} \subseteq \mathcal{V}^{\Sigma}(\mathcal{B})$. This is due to the fact that our aim is to understand whether there is loss. It is clear that when \mathcal{E} is a strict subset of $\mathcal{V}^{\Sigma}(\mathcal{B})$ then loss may occur, but this has nothing to do with the “quality” of the views.

It is now of interest to consider the cases in which some or all of the above inclusions are actually equalities, since these correspond to the notions studied in this paper.

- (1) If $R_{Q,\mathcal{V}}^{\max}$ is exact, i.e., is equivalent to Q (modulo the view definitions), then all three inclusions are actually equalities. Hence, not only do we have losslessness with respect to rewriting but we also have both that the views are lossless with respect to answering and that $R_{Q,\mathcal{V}}^{\max}$ is perfect. Thus exactness of the maximal rewriting is the strongest notion, combining both losslessness of the views and perfectness of the rewriting.
- (2) If $R_{Q,\mathcal{V}}^{\max}$ is perfect, i.e., is equivalent to $\text{cert}_{Q,\mathcal{V}}$, then inclusions (1) and (2) are actually equalities. In this case, we also get that $\text{cert}_{Q,\mathcal{V}}$ has to coincide with $\text{clin}_{Q,\mathcal{V}}$. By Corollary 6.5 we can conclude that the certain answers are expressible as a 2RPQ over \mathcal{B} .
- (3) If \mathcal{V} is lossless with respect to Q , i.e., we have losslessness with respect to answering, then inclusion (3) is actually an equality. Moreover, in this case, since Q is itself a 2RPQ, and hence is linear, then $\text{cert}_{Q,\mathcal{V}}$ has also to be linear and has to coincide with $\text{clin}_{Q,\mathcal{V}}$. Hence inclusion (2) is also an equality. In this case we know that there is not loss of information related to the fact that we are answering the query based on a set of views.
- (4) Finally, if \mathcal{V} is lossy with respect to Q , i.e., we have lossiness with respect to answering, we can check whether inclusion (2) is actually an equality, i.e., whether the certain answers are actually expressible as a 2RPQ over the database. If this is the case, we directly get a characterization of the certain answers in the same language used for expressing the query, namely 2RPQs over the database, and thus in terms that are understandable to the user.

More generally, if \mathcal{V} is lossy with respect to Q and inclusion (2) is a proper inclusion, we would like to provide an explanation for the answers that are actually returned or, equivalently, for the loss of information. Indeed, in this case, we know that there will be at least one view extension such that, in order to show that a tuple is not a certain answer, we need to resort to a non-linear database. It remains to be investigated whether the techniques we provide for doing the check allow one also to extract such a counterexample database to exhibit to the user.

Acknowledgements

Diego Calvanese, Giuseppe De Giacomo, and Maurizio Lenzerini were supported in part by the EU funded FET Project Thinking ONtologiES (TONES), under contract FP6-7603. Giuseppe De Giacomo and Maurizio Lenzerini were further supported by project HYPER, funded by IBM through a Shared University Research (SUR) Award grant, and by MIUR FIRB 2005 project “Tecnologie Orientate alla Conoscenza per Aggregazioni di Imprese in Internet” (TOCALIT). Moshe Y. Vardi was supported in part by NSF grants CCR-9988322, CCR-0124077, CCR-0311326, IIS-9908435, IIS-9978135, EIA-0086264, and ANI-0216467, by US-Israel BSF grant 9800096, by Texas ATP grant 003604-0058-2003, and by a grant from the Intel Corporation.

References

- [1] A.Y. Levy, A.O. Mendelzon, Y. Sagiv, D. Srivastava, Answering queries using views, in: Proc. of the 14th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS’95, 1995, pp. 95–104.
- [2] J.D. Ullman, Information integration using logical views, in: Proc. of the 6th Int. Conf. on Database Theory, ICDT’97, in: Lecture Notes in Computer Science, vol. 1186, Springer, 1997, pp. 19–40.
- [3] S. Abiteboul, O. Duschka, Complexity of answering queries using materialized views, in: Proc. of the 17th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS’98, 1998, pp. 254–265.
- [4] A.Y. Halevy, Answering queries using views: A survey, Very Large Database J. 10 (4) (2001) 270–294.
- [5] O.M. Duschka, M.R. Genesereth, Answering recursive queries using views, in: Proc. of the 16th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS’97, 1997, pp. 109–116.

- [6] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, Answering regular path queries using views, in: Proc. of the 16th IEEE Int. Conf. on Data Engineering, ICDE 2000, 2000, pp. 389–398.
- [7] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, View-based query containment, in: Proc. of the 22nd ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2003, 2003, pp. 56–67.
- [8] F.N. Afrati, C. Li, P. Mitra, Answering queries using views with arithmetic comparisons, in: Proc. of the 21st ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2002, 2002, pp. 209–220.
- [9] S. Flesca, S. Greco, Rewriting queries using views, IEEE Trans. Knowl. Data Eng. 13 (6) (2001) 980–995.
- [10] R. Chirkova, A.Y. Halevy, D. Suciu, A formal perspective on the view selection problem, in: Proc. of the 27th Int. Conf. on Very Large Data Bases, VLDB 2001, 2001, pp. 59–68.
- [11] C. Li, M. Bawa, J.D. Ullman, Minimizing view sets without losing query-answering power, in: Proc. of the 8th Int. Conf. on Database Theory, ICDT 2001, 2001, pp. 99–113.
- [12] S. Grumbach, L. Tininini, On the content of materialized aggregate views, in: Proc. of the 19th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2000, 2000, pp. 47–57.
- [13] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, Lossless regular views, in: Proc. of the 21st ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2002, 2002, pp. 58–66.
- [14] P. Buneman, Semistructured data, in: Proc. of the 16th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS'97, 1997, pp. 117–121.
- [15] S. Abiteboul, Querying semi-structured data, in: Proc. of the 6th Int. Conf. on Database Theory, ICDT'97, 1997, pp. 1–18.
- [16] D. Florescu, A. Levy, A. Mendelzon, Database techniques for the world-wide web: A survey, SIGMOD Record 27 (3) (1998) 59–74.
- [17] S. Abiteboul, P. Buneman, D. Suciu, Data on the Web: From Relations to Semistructured Data and XML, Morgan Kaufmann, 2000.
- [18] P. Buneman, S. Davidson, G. Hillebrand, D. Suciu, A query language and optimization technique for unstructured data, in: Proc. of the ACM SIGMOD Int. Conf. on Management of Data, 1996, pp. 505–516.
- [19] S. Abiteboul, D. Quass, J. McHugh, J. Widom, J.L. Wiener, The Lorel query language for semistructured data, Int. J. Digit. Libr. 1 (1) (1997) 68–88.
- [20] J. Clark, S. DeRose, XML Path Language (XPath) version 1.0—W3C recommendation 16 november 1999, Tech. rep., World Wide Web Consortium, available at <http://www.w3.org/TR/1999/REC-xpath-19991116>, 1999.
- [21] T. Milo, D. Suciu, Index structures for path expressions, in: Proc. of the 7th Int. Conf. on Database Theory, ICDT'99, in: Lecture Notes in Computer Science, vol. 1540, Springer, 1999, pp. 277–295.
- [22] M.F. Fernandez, D. Florescu, J. Kang, A.Y. Levy, D. Suciu, Catching the boat with Strudel: Experiences with a web-site management system, in: Proc. of the ACM SIGMOD Int. Conf. on Management of Data, 1998, pp. 414–425.
- [23] A. Deutsch, M.F. Fernandez, D. Florescu, A. Levy, D. Suciu, XML-QL: A query language for XML, Submission to the World Wide Web Consortium, available at <http://www.w3.org/TR/NOTE-xml-ql>, August 1998.
- [24] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, Query processing using views for regular path queries with inverse, in: Proc. of the 19th ACM SIGACT SIGMOD SIGART Symp. on Principles of Database Systems, PODS 2000, 2000, pp. 58–66.
- [25] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, Reasoning on regular path queries, SIGMOD Rec. 32 (4) (2003) 83–92.
- [26] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, Rewriting of regular expressions and regular path queries, J. Comput. Syst. Sci. 64 (3) (2002) 443–465.
- [27] R. Reiter, On closed world data bases, in: H. Gallaire, J. Minker (Eds.), Logic and Databases, Plenum Publ. Co., 1978, pp. 119–140.
- [28] T. Feder, M.Y. Vardi, The computational structure of monotone monadic SNP and constraint satisfaction, SIAM J Comput. 28 (1999) 57–104.
- [29] G. Grahne, A.O. Mendelzon, Tableau techniques for querying information sources through global schemas, in: Proc. of the 7th Int. Conf. on Database Theory, ICDT'99, in: Lecture Notes in Computer Science, vol. 1540, Springer, 1999, pp. 332–347.
- [30] J.E. Hopcroft, J.D. Ullman, Introduction to Automata Theory, Languages, and Computation, Addison Wesley Publ. Co., 1979.
- [31] D. Calvanese, G. De Giacomo, M. Lenzerini, M.Y. Vardi, View-based query processing and constraint satisfaction, in: Proc. of the 15th IEEE Symp. on Logic in Computer Science, LICS 2000, 2000, pp. 361–371.
- [32] M.Y. Vardi, A temporal fixpoint calculus, in: Proc. of the 15th ACM SIGPLAN-SIGACT Symp. on Principles of Programming Languages, POPL'88, San Diego, CA, USA, 1988, pp. 250–259.
- [33] M.Y. Vardi, A note on the reduction of two-way automata to one-way automata, Inform. Process. Lett. 30 (5) (1989) 261–264.